# Parameter estimation of hidden Markov models: comparison of EM and quasi-Newton methods with a new hybrid algorithm

Sidonie Foulon, Thérèse Truong, Anne-Louise Leutenegger & Hervé Perdry

Sidonie Foulon - MaSeMo workshop - July, 4th, 2025







CESP

# Plan

- 1. Introduction
- 2. Methods
  - a. Algorithms
  - b. Examples / data
- 3. Results
  - a. Application on the examples
  - b. Convergence points
- 4. Conclusion & Discussion

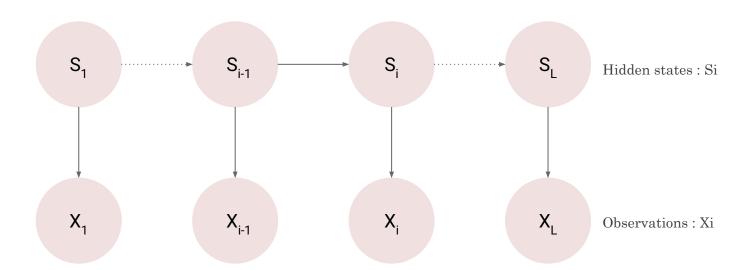
## Introduction: Hidden Markov models

A characteristic of a Markov chain  $S_i$  is the homogeneity of the sequence.  $\Rightarrow$  sometimes contradicted by the observations.

### To model such situations:

- \* add an **observable** layer of variables to the model
- \* related to the sequence of **hidden** variables

# $\Rightarrow$ Hidden Markov Models (HMM)



# Introduction: Hidden Markov models

A characteristic of a Markov chain  $S_i$  is the homogeneity of the sequence.

 $\Rightarrow$  sometimes contradicted by the observations.

To model such situations:

- \* add an **observable** layer of variables to the model
- \* related to the sequence of **hidden** variables

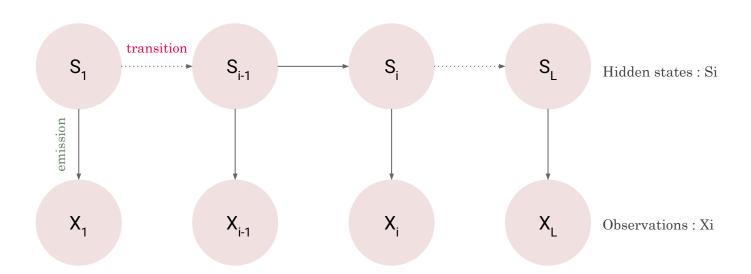
Transition probabilities:

$$T(s,t) = \mathbb{P}(S_i = t | S_{i-1} = s)$$

Emission probabilities:

$$E(x,s) = \mathbb{P}(X_i = x | S_i = s)$$

 $\Rightarrow$  Hidden Markov Models (HMM)



# Aim: HMM parameters estimation

# **Introduction: Parameter estimation**

Classical optimisation methods to estimate HMM parameters  $\theta$ 

### Expectation-Maximisation (EM) algorithm:

Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

Estimation of parameters by a variant of EM algorithm: Baum-Welch algorithm.

### <u>Direct maximisation of likelihood (DML)</u>:

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

Likelihood maximisation by a quasi-Newton type method (L-BFGS-B).

# Introduction: Parameter estimation

# Classical optimisation methods to estimate HMM parameters $\theta$

### Expectation-Maximisation (EM) algorithm:

Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

Estimation of parameters by a variant of EM algorithm: Baum-Welch algorithm.

### <u>Direct maximisation of likelihood (DML)</u>:

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

Likelihood maximisation by a quasi-Newton type method (L-BFGS-B).

### Hybrid method:

Creation of QNEM, a mix of EM algorithm and a quasi-Newton type method (BFGS).

# **Introduction: Parameter estimation**

Classical optimisation methods to estimate HMM parameters  $\theta$ 

### Expectation-Maximisation (EM) algorithm:

Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

Estimation of parameters by a variant of EM algorithm: Baum-Welch algorithm.

### <u>Direct maximisation of likelihood (DML)</u>:

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

Likelihood maximisation by a quasi-Newton type method (L-BFGS-B).

## EM acceleration:

Varadhan and Roland, 2008

SQUAREM method (squared iterative methods).

### Hybrid method:

Creation of QNEM, a mix of EM algorithm and a quasi-Newton type method (BFGS).

⇒ Comparison of the 4 methods

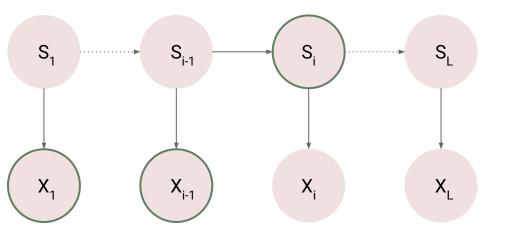
Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- ❖ 2 possible forward algorithms:
  - > joint probabilities

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

 $a_i(s) = \mathbb{P}(S_i = s, X_1^{i-1} = x_1^{i-1})$ 

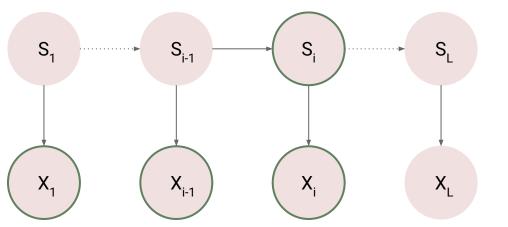
- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory
    Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities



Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

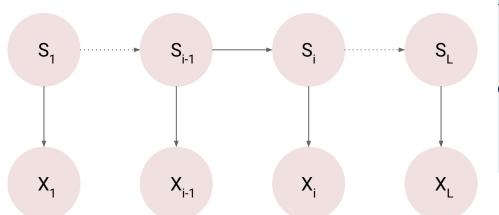
- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities

$$a_i(s) = \mathbb{P}(S_i = s, X_1^{i-1} = x_1^{i-1})$$
  
 $b_i(s) = \mathbb{P}_{\theta}(S_i = s, X_1^i = x_1^i)$ 



Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities



$$a_i(s) = \mathbb{P}(S_i = s, X_1^{i-1} = x_1^{i-1})$$
  
 $b_i(s) = \mathbb{P}_{\theta}(S_i = s, X_1^i = x_1^i)$ 

### 1. Initialisation:

$$a_1(s) = \pi(s) \tag{1}$$

al distribution of state s. with  $\pi(s)$  the init

2. Recursion:

$$\sum_{t} b_{i-1}(t) T_{\theta}(t, s) \tag{3}$$

$$= a_i(s)E_{\theta}(x_i, s) \tag{4}$$

here the sum in (3) is over t.

d can be computed as

$$=x_1^L)=\sum a_L(s)$$
 . (5)

$$a_i(s) =$$

$$b_i(s)$$
 :

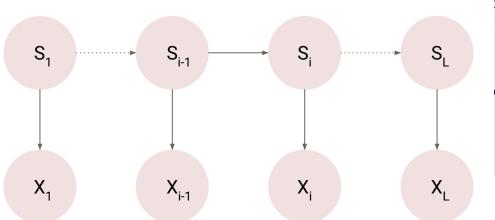
for  $i \in 2, ..., L$ , w all possible states

Finally, the likelihoo

$$L(\theta; \mathbf{X}) = \mathbb{P}_{\theta}(X_1^I)$$

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian matrix at each iteration
  - ➤ method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities
    - need to compute everything on log scale



$$a_i(s) = \mathbb{P}(S_i = s, X_1^{i-1} = x_1^{i-1})$$
  
 $b_i(s) = \mathbb{P}_{\theta}(S_i = s, X_1^i = x_1^i)$ 

### 1. Initialisation:

$$a_1(s) = \pi(s) \tag{1}$$

al distribution of state s. with  $\pi(s)$  the init

2. Recursion:

$$\sum_{t} b_{i-1}(t) T_{\theta}(t, s) \tag{3}$$

$$= a_i(s)E_{\theta}(x_i, s) \tag{4}$$

here the sum in (3) is over t.

d can be computed as

$$=x_1^L) = \sum a_L(s)$$
 . (5)

$$a_i(s) =$$

$$b_i(s)$$
 :

for  $i \in 2, ..., L$ , w all possible states

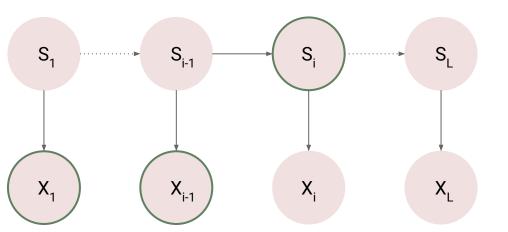
Finally, the likelihoo

$$L(\theta; \mathbf{X}) = \mathbb{P}_{\theta}(X_1^I)$$

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

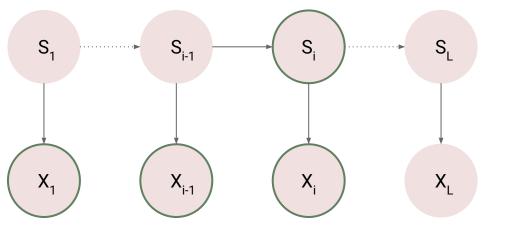
 $\alpha_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^{i-1} = x_1^{i-1})$ 

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities
    - need to compute everything on log scale
  - > conditional probabilities



Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory
    Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities
    - need to compute everything on log scale
  - > conditional probabilities

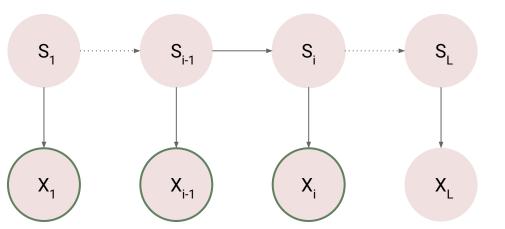


$$\alpha_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^{i-1} = x_1^{i-1})$$

$$\beta_i(s) = \mathbb{P}_{\theta}(S_i = s|X_1^i = x_1^i)$$

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

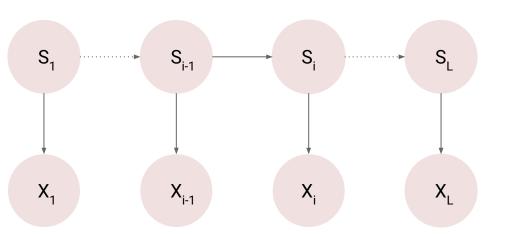
- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - > method used: L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - > joint probabilities
    - need to compute everything on log scale
  - > conditional probabilities



$$\alpha_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^{i-1} = x_1^{i-1})$$
$$\beta_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^i = x_1^i)$$
$$\gamma_i = \log \left(\mathbb{P}_{\theta}(X_1^i = x_1^i)\right).$$

Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970

- Quasi-Newton methods
  - maximisation of differentiable scalar functions
  - ➤ avoiding computing the Hessian matrix
     ⇒ approximating the inverse Hessian
     matrix at each iteration
  - ➤ method used : L-BFGS-B (Limited-memory Broyden-Fletcher-Goldfarb-Shanno with Bound constraints)
- ❖ Forward step → likelihood
  - → direct maximisation with L-BFGS-B
- 2 possible forward algorithms:
  - joint probabilities
    - need to compute everything on log scale
  - > conditional probabilities



$$\alpha_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^{i-1} = x_1^{i-1})$$

$$\beta_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^i = x_1^i)$$

$$\gamma_i = \log \left( \mathbb{P}_{\theta}(X_1^i = x_1^i) \right).$$

1. Initialisation:

$$\alpha_1(s) = \pi_s \tag{8}$$

$$\beta_1(s) = \frac{E_{\theta}(x_1, s) \times \alpha_1(s)}{\sum_i E_{\theta}(x_1, i) \times \alpha_1(i)}$$
(9)

$$\gamma_1 = \log \left( \sum_s \alpha_1(s) E_{\theta}(x_1, s) \right).$$
 (10)

2. Recursion:

$$\alpha_i(s) = \sum_t T_{\theta}(t, s) \times \beta_{i-1}(t) \qquad (11)$$

$$\beta_i(s) = \frac{E_{\theta}(x_i, s) \times \alpha_i(s)}{\sum_t E_{\theta}(x_i, t) \times \alpha_i(t)}$$
(12)

$$\gamma_i = \gamma_{i-1} + \log \left( \sum_s E_{\theta}(x_i, s) \alpha_i(s) \right)$$
 (13)

for  $i \in 2, ..., L$ 

3. Log-likelihood :  $\ell(\theta; \mathbf{X}) = \gamma_L$ 

Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

**E step**: knowing  $\theta$ , compute the probabilities of the hidden states along the chain given all the observations  $\Rightarrow$  forward-backward algorithm

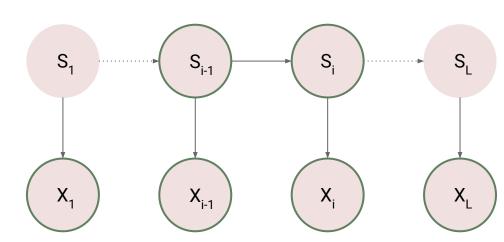
- conditional forward algorithm
- backward algorithm

Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

**E** step: knowing  $\theta$ , compute the probabilities of the hidden states along the chain given all the observations  $\Rightarrow$  forward-backward algorithm

- conditional forward algorithm
- backward algorithm:

$$\delta_i(s,t) = \mathbb{P}_{\theta}(S_{i-1} = s, S_i = t | X_1^L = x_1^L)$$



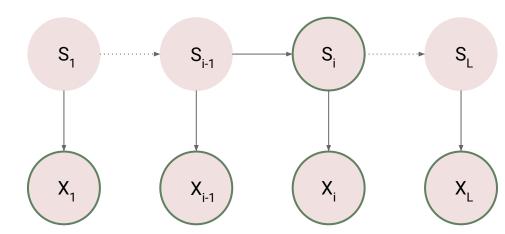
Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

**E** step: knowing  $\theta$ , compute the probabilities of the hidden states along the chain given all the observations  $\Rightarrow$  forward-backward algorithm

- conditional forward algorithm
- backward algorithm:

$$\delta_i(s,t) = \mathbb{P}_{\theta}(S_{i-1} = s, S_i = t | X_1^L = x_1^L)$$

$$\varphi_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^L = x_1^L)$$



Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

**E** step: knowing  $\theta$ , compute the probabilities of the hidden states along the chain given all the observations  $\Rightarrow$  forward-backward algorithm

- conditional forward algorithm
- backward algorithm :

$$\delta_i(s,t) = \mathbb{P}_{\theta}(S_{i-1} = s, S_i = t | X_1^L = x_1^L)$$

$$\varphi_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^L = x_1^L)$$

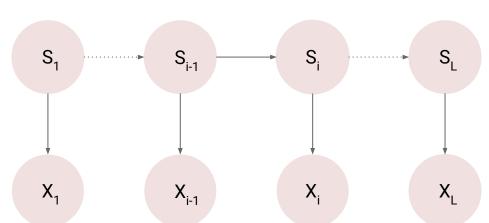
1. Initialisation:

$$\varphi_L(s) = \beta_L(s)$$
 (16)

2. Recursion:

$$\delta_i(s,t) = T_{\theta}(s,t) \times \frac{\beta_{i-1}(s)}{\alpha_i(t)} \times \varphi_i(t)$$
 (17)

$$\varphi_{i-1}(s) = \sum_{t} \delta_i(s, t) \tag{18}$$



Dempster et al., 1977; Baum and Petrie, 1966; Welch, 2003

**E** step: knowing  $\theta$ , compute the probabilities of the hidden states along the chain given all the observations  $\Rightarrow$  forward-backward algorithm

- conditional forward algorithm
- backward algorithm :

$$\delta_i(s,t) = \mathbb{P}_{\theta}(S_{i-1} = s, S_i = t | X_1^L = x_1^L)$$
$$\varphi_i(s) = \mathbb{P}_{\theta}(S_i = s | X_1^L = x_1^L)$$

1. Initialisation:

$$\varphi_L(s) = \beta_L(s) \tag{16}$$

2. Recursion:

$$\delta_i(s,t) = T_{\theta}(s,t) \times \frac{\beta_{i-1}(s)}{\alpha_i(t)} \times \varphi_i(t)$$
 (17)

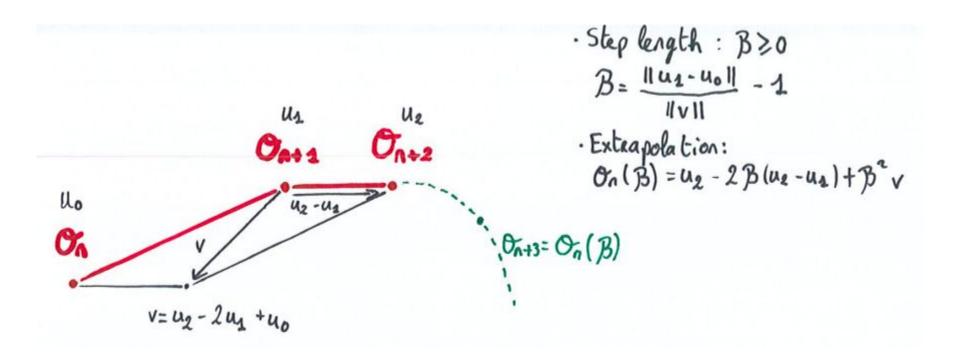
$$\varphi_{i-1}(s) = \sum_{t} \delta_i(s, t) \tag{18}$$

**M** step: re-estimate  $\theta$  using these probabilities

- for each iteration of EM algorithm, estimation of the parameter θ for the next iteration
- > if we knew the sequence of hidden state, we could compute  $\theta$  with counting
- > but the hidden states are unknown  $\Rightarrow$  estimation of θ with probabilities

# SQUAREM Varadhan and Roland, 2008

- ❖ SQUAREM = Squared iterative Methods
- EM acceleration algorithm
- Simple and Globally Convergent Methods for Accelerating the Convergence of Any EM Algorithm, Varadhan et al., 2008, Scand. J. Stat.
- ❖ Aim: maintain the stability of EM while accelerating its convergence speed



# Method: BFGS

## Algorithm 2 BFGS algorithm

- 0. Input:  $\theta^{(0)} = \text{starting point.}$ Set k = 0 and  $H_0 = \mathbb{I}$ .
- 1. Let  $p_k = -H_k \cdot \nabla f_k$  the search direction. Perform a line search for  $\alpha$  such that the Wolfe conditions are satisfied

$$f(\theta_k + \alpha \cdot p_k) \le f(\theta_k) + c_1 \cdot \alpha(\nabla f_k)' p_k ;$$
  
$$\nabla f(\theta_k + \alpha \cdot p_k)' p_k \ge c_2 (\nabla f_k)' p_k .$$

- 2. Let  $\theta^{(k+1)} = \theta^{(k)} + \alpha \cdot p_k$ ,  $s_k = \theta^{(k+1)} \theta^{(k)}$  and  $y_k = \nabla f_{k+1} \nabla f_k$ .
- 3. BFGS (quasi-Newton) step: Update  $H_k$  using the BFGS formula

$$H_{k+1} = H_k - \frac{H_k \cdot y_k \cdot y_k' \cdot H_k}{y_k' \cdot H_k \cdot y_k} + \frac{s_k \cdot s_k'}{y_k' \cdot s_k}.$$

Let k = k + 1. Repeat 1. to 3. until convergence.

**Curvature condition :** 

$$s_k' \cdot y_k > 0$$

with 
$$s_k = \theta_{k+1} - \theta_k$$
  
and  $y_k = \nabla f_{k+1} - \nabla f_k$ 

- With a complex line search, this condition is always met
- If the curvature condition is met,
   H<sub>k</sub> stay positive definite

# Proposed method: QNEM

### Algorithm 1 QNEM algorithm

- 0. Input:  $\theta^{(0)} = \text{starting point.}$ Set k = 0 and  $H_0 = \mathbb{I}$ .
- 1. Baum-Welch (EM) step: Compute

$$\theta^{(k+1)} = BW(\theta^{(k)}) .$$

- 2. Let  $s_k = \theta^{(k+1)} \theta^{(k)}$  and  $y_k = \nabla f_{k+1} \nabla f_k$ . Repeat 1 & 2 until the curvature condition  $s'_k \cdot y_k > 0$  is satisfied, or the convergence criterion (Section 3.5) is met.
- 3. BFGS (quasi-Newton) step: Update  $H_k$  using the BFGS formula

$$H_{k+1} = H_k - \frac{H_k \cdot y_k \cdot y_k' \cdot H_k}{y_k' \cdot H_k \cdot y_k} + \frac{s_k \cdot s_k'}{y_k' \cdot s_k}.$$

Let k = k + 1.

4. Let  $p_k = -H_k \cdot \nabla f_k$  the search direction. Perform a backtracking search for  $\alpha$  such that the Armijo condition is satisfied

$$f(\theta_k + \alpha \cdot p_k) \le f(\theta_k) + c \cdot \alpha(\nabla f_k)' p_k$$
.

5. Let  $\theta^{(k+1)} = \theta^{(k)} + \alpha \cdot p_k$ . If the curvature condition is satisfied, repeat 3 to 5 until the convergence criterion is met. Else, reinitialise  $H_k = \mathbb{I}$  and go back to 1.

- Simple line search = backtracking (Armijo condition)
- If the curvature condition is not
   met ⇒ switch to EM

# **Example overview**

## <u>Umbrella example:</u>

Toy example: weather forecast based on the presence of an umbrella

- $\rightarrow$  2 hidden states
- → discrete observations (2 values)
- $\rightarrow$  N = 56

### Geyser example:

Example based on open access data "Old Faithful Geyser": chain on the duration of geyser eruptions

$$\rightarrow$$
 N = 272

### 2 sub-examples:

- dichotomised times :
  - $\rightarrow$  3 hidden states
  - → discrete observations (2 values)
- continuous times:
  - $\rightarrow$  3 hidden states
  - → continuous observations

### Genetics example:

Motivating problem for this work (part of my PhD): identification of homozygous segments. Data simulated from 1006 european HGDP-CEPH (Human Genome Diversity Project) haplotypes.

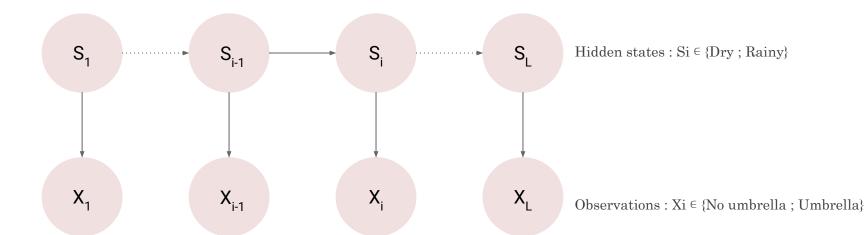
- $\rightarrow$  2 hidden states
- → discrete observations (3 values)
- $\rightarrow$  N = 1050

# <u>Umbrella example:</u>

Toy example: weather forecast based on the presence of an umbrella

- $\rightarrow$  2 hidden states
- → discrete observations (2 values)
- $\rightarrow$  N = 56

Transition matrix	$S_i = D$	$S_i = R$
$S_{i-1} = D$	1-a	a
$S_{i-1} = \mathbf{R}$	$\mathbf{a}$	1-a
Emission matrix	$S_i = D$	$S_i = R$
$X_i = N$	1-b	b
$X_i = U$	b	1-b



# Results: Umbrella example

Umbrella example		Nb of iterations						Mean nb of steps		Percent of
	Min	Q1	Q2	Mean	Q3	Max	Fw	Bw	time (s)	convergence
Quasi-Newton	4	13	16	15.77	18	37	15.77	0	0.03	100
Baum-Welch	3	18	24	27.95	31	338	27.95	27.95	0.06	100
SQUAREM	3	20	24	26.76	29	179	17.70	17.67	0.04	100
QNEM	2	8	11	11.32	13	118	12.63	2.80	0.03	100

The mean running times are in seconds.

Backward: ~1.7ms

Forward Baum-Welch & SQUAREM : ~0.5ms

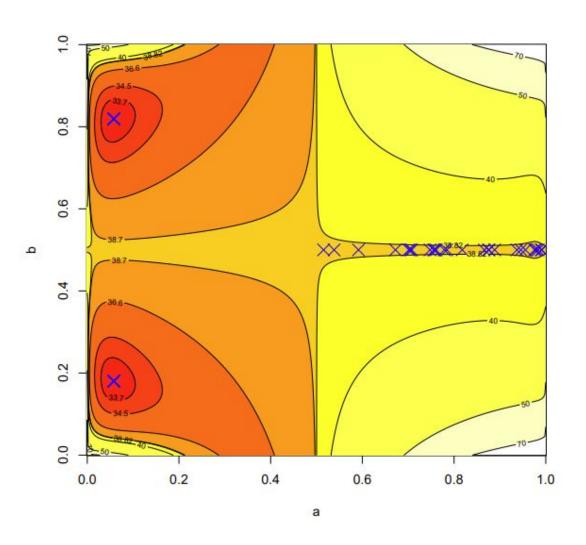
Forward quasi-Newton & QNEM : ~1.1ms

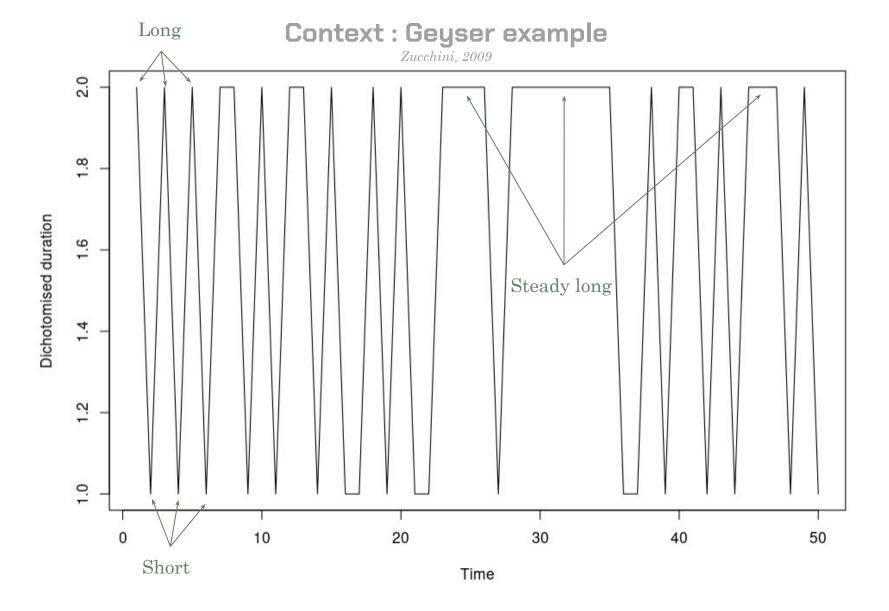
# Results: Umbrella example

Likelihood	Parar	neters	Percent of permutations								
	a	b	Quasi-Newton	Baum-Welch	SQUAREM	QNEM					
33.3	0.06	$0.18 \\ 0.82$	89.9	70.1	69.9	70.1					
33.8	ND	0.5	10.1	29.9	30.1	29.3					
Other			0	0	0	0.6					

Table 8: This table shows the different convergence points reached with the umbrella example and their percent. ND = Non determined.

# Log-likelihood with 100 convergence points





Plot of dichotomised geyser eruptions' durations along time.

- "1" code for eruption shorter than 3 minutes.
- "2" code for eruption longer than 3 minutes.

### Geyser example:

Example based on open access data "Old Faithful Geyser": chain on the duration of geyser eruptions

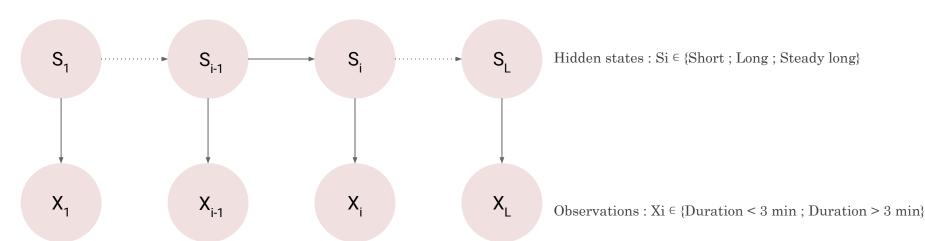
$$\rightarrow$$
 N = 272

# 2 sub-examples:

- dichotomised times:
  - $\rightarrow$  3 hidden states
  - → discrete observations (2 values)
- continuous times :
  - $\rightarrow$  3 hidden states
  - $\rightarrow$  continuous observations

Transition matrix	$S_i = S$	$S_i = L$	$S_i = Sl$
$S_{i-1} = S$	0	1-a	a
$S_{i-1} = L$	1	0	0
$S_{i-1} = \operatorname{Sl}$	1-b	0	b

Emission matrix	$S_i = S$	$S_i = L$	$S_i = \operatorname{Sl}$
$X_i = \text{Dinf3}$	1-c	1-d	1-e
$X_i = \text{Dsup3}$	$\mathbf{c}$	d	e



# Results: Dichotomised geyser example

Geyser example		1	Nb of	iteration	ns		Mean nb of steps		Mean	Percent of
(dichotomised)	Min	Q1	Q2	Mean	Q3	Max	Fw	Bw	time (s)	convergence
Quasi-Newton	12	18	21	24.23	26	66	24.23	0	0.33	100
Baum-Welch	12	82	139	126	161	324	126	126	1.21	100
SQUAREM	15	71	113	108	140	290	78.59	71.89	0.71	100
QNEM	6	14	16	16.27	19	64	17.55	1.62	0.26	99.8

The mean running times are in seconds.

Backward:~8ms

Forward Baum-Welch & SQUAREM: ~1.5ms

Forward quasi-Newton & QNEM : ~11.5ms

# Results: Dichotomised geyser example

Likelihood			Paramete	ers		Percent of permutations				
a b c d	e	Quasi-Newton	Baum-Welch	SQUAREM	QNEM					
144.5	0.79	0.57	0*	1*	0.95	52.6	50.6	50.8	42.6	
149.5	0.58	0.50	1*	0*	0.57	41.2	42.4	42.4	36.2	
172.8	0* ND	ND 0*	0.55 / 0.74	0.74 / 0.55 ND	ND	6.2	7.0	0.9 5.9	4.4 4.5	
Other						0	0	0	11.1	

Table 10: This table shows the different convergence points reached with the dichotomised geyser example and their percent. ND = Non determined. \* = These values were rounded for Quasi-Newton results because of the bound-constraints.

### Geyser example:

Example based on open access data "Old Faithful Geyser": chain on the duration of geyser eruptions

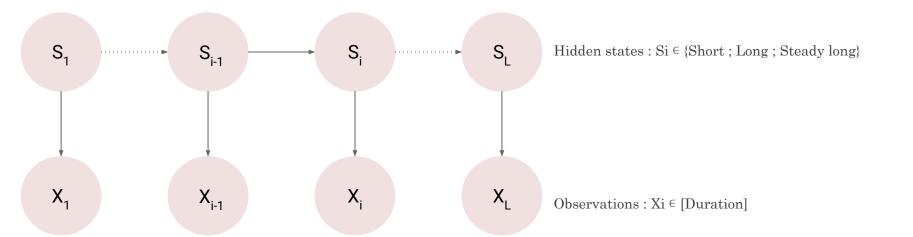
$$\rightarrow$$
 N = 272

# 2 sub-examples:

- dichotomised times :
  - $\rightarrow$  3 hidden states
  - → discrete observations (2 values)
- continuous times:
  - $\rightarrow$  3 hidden states
  - $\rightarrow$  continuous observations

Transition matrix	$S_i = S$	$S_i = L$	$S_i = Sl$
$S_{i-1} = S$	0	1-a	a
$S_{i-1} = L$	1	0	0
$S_{i-1} = \operatorname{Sl}$	1-b	0	b

The emission densities are the Gaussian densities with parameters  $\mu_x$  and  $\sigma_x$  for state x.



# Results: Continuous geyser example

Geyser example		Nb of iterations						Mean nb of steps		Percent of
(continuous)	Min	Q1	Q2	Mean	Q3	Max	Fw	Bw	time (s)	convergence
Quasi-Newton	21	60	78	79.41	98	159	79.41	0	1.47	93.2
Baum-Welch	8	21	23	25.32	26	129	25.32	25.32	0.24	100
SQUAREM	12	30	32	34.04	35	118	23.23	22.59	0.22	100
QNEM	4	19	23	24.01	27	87	31.61	1.89	0.66	99.5

The mean running times are in seconds.

Backward:~8ms

Forward Baum-Welch & SQUAREM : ~1.5ms

Forward quasi-Newton & QNEM : ~17.4ms

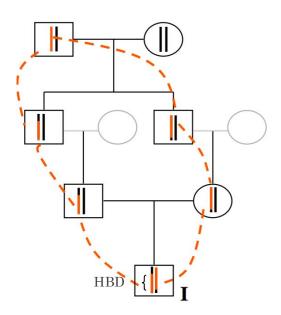
# Results: Continuous geyser example

Likelihood	Parameters								Percent of permutations				
	a	b	$\mu_s$	$\mu_l$	$\mu_{sl}$	$\sigma_s$	$\sigma_l$	$\sigma_{sl}$	Quasi-Newton	Baum-Welch	SQUAREM	QNEM	
265.7	0.61	0.65	2.0	4.58	4.09	0.22	0.24	0.64	19.1	40.7	40.5	35.7	
275.1	0.54	0.46	4.44	1.98	3.23	0.3	0.19	1.03	29.2	38.6	38.6	30.6	
303.2	0.93	0.34	2.37	1.91	4.34	0.77	0.20	0.37	21.9	9.0	9.6	9.9	
316.2	1*	0.28	4.45	ND	2.79	0.30	ND	1	0.1	0	0	6.4	
335	$0.44 \\ 0.56$	0*	3.72	4.24 2.0	$\frac{2.0}{4.24}$	1.05	$0.45 \\ 0.22$	$0.22 \\ 0.45$	7.5	5.5	4.9	5.2	
Other									22.2	6.2	6.4	12.2	

Table 12: This table shows the different convergence points reached with the continuous geyser example and their percent. ND = Non determined. \* = These values were rounded for Quasi-Newton results because of the bound-constraints.

### **Context: Genetics example**

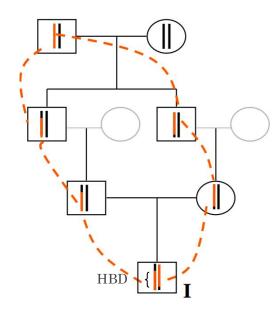
- ❖ Identification of rare recessive variants involved in multifactorial disease
- Analysis on consanguineous individuals (offspring of relatives), which are more likely to carry this type of variants
- **❖ Inbreeding coefficient** *f*: probability that 2 alleles at a locus drawn in an individual's genome are identical and come from a common ancestor to his parents



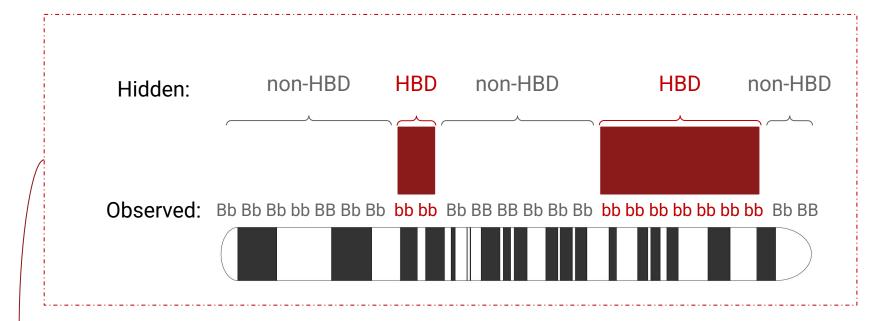
Genealogy of an offspring of first cousins

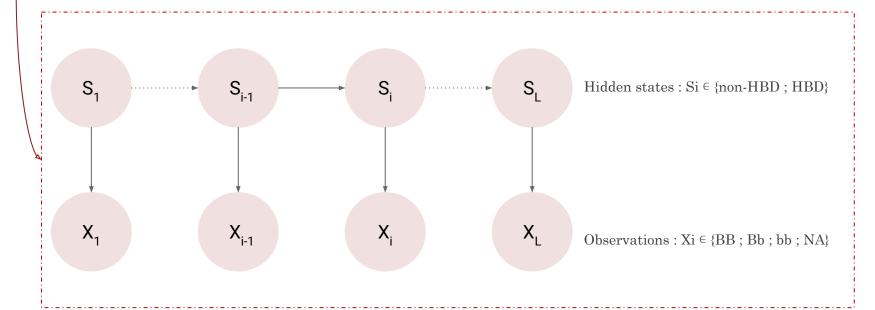
## **Context: Genetics example**

- In absence of family data: identified through his genome which carries long homozygous segments: HBD segments
- \* We introduce a as: the mean length of HBD segments is 1/(a(1-f)) and the mean length of non-HBD segments is 1/(af) (in cM)
- ❖ HBD segments inference ⇒ Hidden Markov model (HMM)



Genealogy of an offspring of first cousins





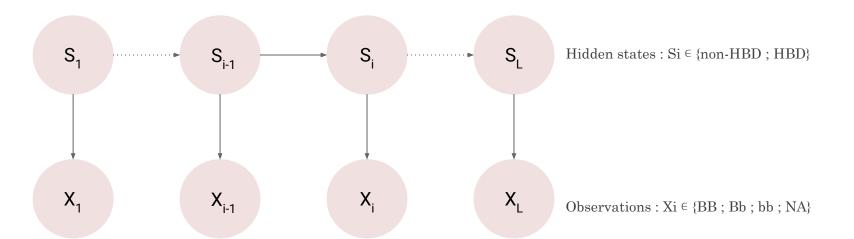
#### Genetics example:

Motivating problem for this work (part of my PhD): identification of homozygous segments. Data simulated from 1006 european HGDP-CEPH (Human Genome Diversity Project) haplotypes.

- $\rightarrow$  2 hidden states
- → discrete observations (3 values)
- $\rightarrow$  N = 1050

Transition matrix	$S_i = \mathrm{nHBD}$	$S_i = \overline{\mathrm{HBD}}$
$S_{i-1} = \text{nHBD}$	$(1 - \exp^{-ad})(1 - f) + \exp^{-ad}$	$(1 - \exp^{-ad})f$
$S_{i-1} = \text{HBD}$	$(1 - \exp^{-ad})(1 - f)$	$(1 - \exp^{-ad})f + \exp^{-ad}$

$S_i = \text{nHBD}$	$S_i = \mathrm{HBD}$
$p_B^2$	$(1 - \epsilon)p_B + \epsilon \cdot p_B^2$
Service Control of the Control of th	$2\epsilon p_B \cdot p_b$
$p_{ar{b}}$	$(1-\epsilon)p_b + \epsilon \cdot p_b^2$
	$\frac{S_i = \text{nHBD}}{p_B^2}$ $2 \cdot p_B \cdot p_b$ $p_b^2$ 1



# Results: HBD segments example

HBD segments		]	Nb of	iteration	ıs		Mean nb of steps		Mean	Percent of	
example	Min	Q1	Q2	Mean	Q3	Max	$F_{\mathbf{W}}$	Bw	time (s)	convergence	
Quasi-Newton	7	11	14	15.36	18	38	15.36	0	0.32	100	
Baum-Welch	5	63	73	68.53	78	86	68.53	68.53	2.40	100	
SQUAREM	5	50	57	56.19	63	88	38.66	38.66	1.37	100	
QNEM	5	9	11	11.99	14	27	12.83	1.40	0.32	100	

The mean running times are in seconds.

Backward: ~30.3ms

Forward Baum-Welch & SQUAREM : ~5.3ms

Forward quasi-Newton & QNEM: ~20.2ms

#### Conclusion & discussion

- Proposed QNEM ⇒ take advantage of Expectation-Maximisation algorithm and Direct maximisation of the likelihood
  - Expectation-Maximisation algorithm : Baum-Welch algorithm
    - + stay close to the solution
    - takes time to converge
  - Direct maximisation of the likelihood : quasi-Newton algorithm
    - + converge faster
    - needs proper initialisation
- Evaluated SQUAREM to measure the extent of the acceleration
- QNEM showed the best results in all example except continuous geyser example where EM / SQUAREM are best
- For genetics example quasi-Newton algorithm and QNEM are equivalent
- No uniformly better algorithm
  - speed of convergence
  - o global vs. local maximum

# **QNEM** availability

Submitted article (under revision) to Journal of Computational and Applied Mathematics

 $\Rightarrow$  preprint available in HAL:

https://hal.science/hal-04685772v2

QNEM implemented in R

 $\Rightarrow$  R package *steveHMM* available on GitHub:

https://github.com/SidonieFoulon/steveHMM



To preprint



To R package

# MERCI!

NeuroDiderot / ICM

CESP Exposome & heredity, genetic group



 $CESP\ HiDiBiostat$ 





# Supplementary

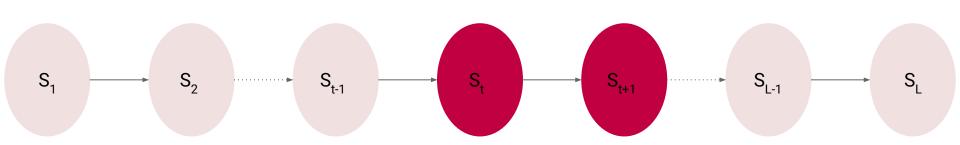
## Introduction: From Markov models ...

A sequence of random variables measured at successive moments  $S_i$  is a **Markov chain** if it satisfies the Markov property:

 $\bigstar$  to predict all the  $S_i$  subsequent to the time t, the information collected for  $i \le t$  is completely included in the single value  $S_t$ 

$$\mathbb{P}(S_{t+1}|S_{1:t}) = \mathbb{P}(S_{t+1}|S_t)$$

- ★ S<sub>i</sub>'s are not independent
- $\star$  S's are independent conditionally on the previous observations

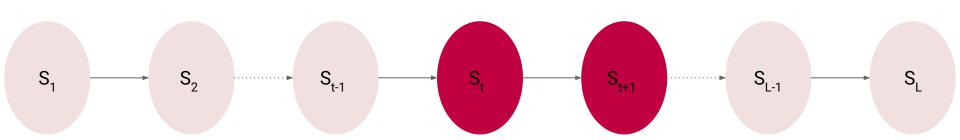


#### Introduction: From Markov models ...

Chain is stationary if : {  $S_{1\text{+d}},\,...,\,S_{\text{t+d}}\}$  follow the same distribution as {  $S_1,\,...,\,S_t\}$ 

Homogeneity of the sequence:

- Markov chain converge fast to this stationary distribution
- **transition probabilities**  $\mathbb{P}(S_{i+1}|S_i)$  are identical in every point of the sequence
- > sometimes contradicted by the observations



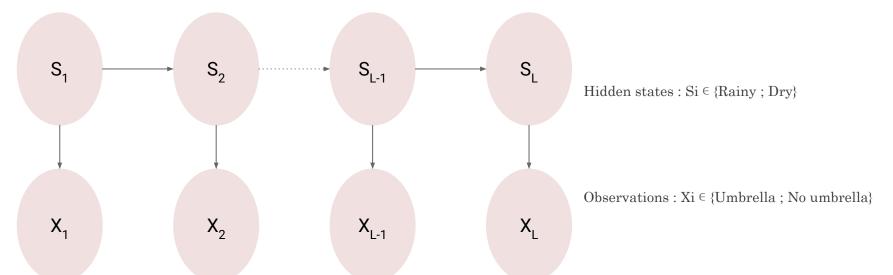
#### Introduction: ... to Hidden Markov models

#### To model such situations:

- \* add an **observable** layer of variables to the model
- \* related to the sequence of **hidden** variables

#### ⇒ Hidden Markov Models (HMM)

- → ex: proportion of homozygous genotypes is almost 1 in HBD segments unlike the rest of the genome
  - ♦ HMM to reconstruct an hidden information (HBD segments) from an observable information (genotypes)



#### Thesis context

- Identification of rare recessive variants involved in multifactorial disease
- Analysis on consanguineous individuals (offspring of relatives), which are more likely to carry this type of variants
- ❖ Inbreeding coefficient f: probability that 2 alleles at a locus drawn in an individual's genome are identical and come from a common ancestor to his parents
- In absence of family data: identified through his genome which carries long homozygous segments: HBD segments
- ❖ We introduce a as: the mean length of HBD segments is 1/(a(1-f)) and the mean length of non-HBD segments is 1/(af) (in cM)
- ♦ HBD segments inference ⇒ Hidden Markov model (HMM)

consanguineous individual Observed: Hidden: non-HBD HBD non-HBD HBD non-HBD

Representation of the HBD segments in the chromosome 8 of a

## Methods: EM algorithm

**E** step: knowing  $\theta$ , compute the probabilities of the hidden states along the chain knowing all the observations

⇒ forward-backward algorithm

- Forward initialize with  $\alpha_i(s) = P_{\theta}(S_1 = s)$ . Then for each position  $i \in \{2,...,L\}$ , compute:
  - $\begin{array}{ll}
    \circ & \alpha_{i}(s) = P_{\theta}(S_{i} = s \mid X_{1}, ..., X_{i-1}) \\
    \circ & \beta_{i}(s) = P_{\alpha}(S_{i} = s \mid X_{1}, ..., X_{i})
    \end{array}$

- Backward initialize with  $\phi_L(s) = \beta_L(s)$ . Then for each position  $i \in \{L-1,...,1\}$ , compute :
  - $\circ \quad \boldsymbol{\delta}_{i}(s,t) = P_{\boldsymbol{\theta}}(S_{i-1} = s, S_{i} = t \mid \mathbf{X})$
  - $\circ \qquad \mathbf{\phi}_{i}(s) = P_{\mathbf{\theta}}(S_{i} = s \mid \mathbf{X})$

**M** step: re-estimate  $\theta$  using these probabilities

The likelihood of  $\theta$  is:

$$\begin{split} L(\pmb{\theta} \ ; \ \mathbf{S} = \mathbf{s}, \ \pmb{X} = \pmb{x}) &= P_{\pmb{\theta}}(S_1 = s_1) \bullet P_{\pmb{\theta}}(S_2 = s_2 \ | \ S_1 = s_1) \bullet \dots \bullet P_{\pmb{\theta}}(S_L = s_L \ | \ S_{L-1} = s_{L-1}) \\ &\bullet P_{\pmb{\theta}}(X_1 = x_1 \ | \ S_1 = s_1) \bullet \dots \bullet P_{\pmb{\theta}}(X_L = x_L \ | \ S_L = s_L) \end{split}$$

$$\begin{array}{l} \boldsymbol{\rightarrow} & l(\boldsymbol{\theta} \; ; \; \mathbf{S} = \mathbf{s}, \; \mathbf{X} = \mathbf{x}) = \log \; P_{\boldsymbol{\theta}}(S_{_{1}} = s_{_{1}}) + \log \; P_{\boldsymbol{\theta}}(S_{_{2}} = s_{_{2}} \; | \; S_{_{1}} = s_{_{1}}) + \ldots + \log \; P_{\boldsymbol{\theta}}(S_{_{L}} = s_{_{L}} \; | \; S_{_{L-1}} = s_{_{L-1}}) \\ & + \log \; P_{\boldsymbol{\theta}}(X_{_{1}} = x_{_{1}} \; | \; S_{_{1}} = s_{_{1}}) + \ldots + \log \; P_{\boldsymbol{\theta}}(X_{_{L}} = x_{_{L}} \; | \; S_{_{L}} = s_{_{L}}) \\ \end{array}$$

Expected value of the log-likelihood:

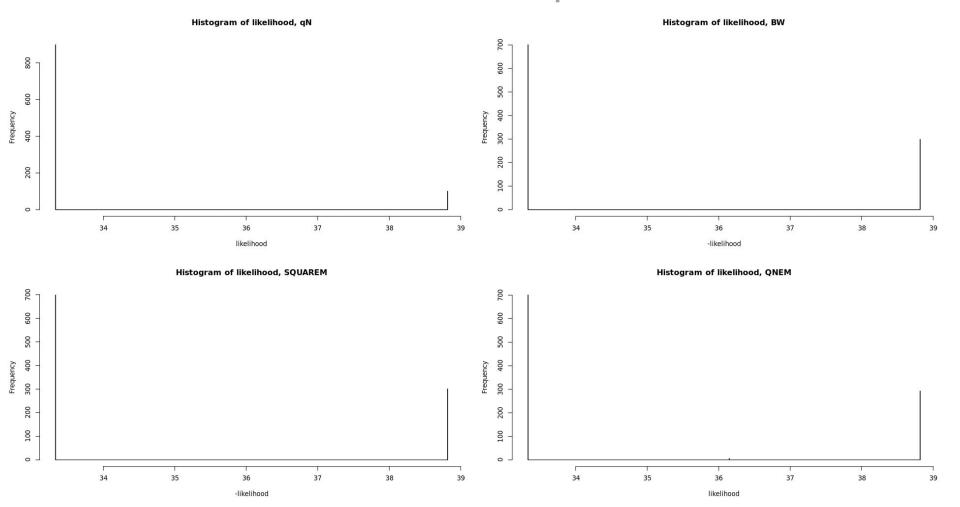
$$\begin{aligned} \mathbf{Q}(\mathbf{\theta}\;;\; \mathbf{\theta}^{(\mathrm{w})}) &= \mathrm{E}(\;\mathrm{l}(\mathbf{\theta}\;;\; \mathbf{S}=\mathbf{s},\, \mathbf{X}=\mathbf{x}) \;\mid\; \mathbf{\theta}^{(\mathrm{w})}) \\ &= \mathrm{\Sigma_{s}} \; \mathrm{log}\; \mathrm{P}_{\mathbf{\theta}}(\mathrm{S}_{1}=\mathbf{s}) \;\bullet\; \mathbf{\phi}_{1}(\mathbf{s}) + \mathrm{\Sigma_{s,t}} \; \mathrm{log}\; \mathrm{P}_{\mathbf{\theta}}(\mathrm{S}_{2}=\mathbf{t}\;\mid\; \mathrm{S}_{1}=\mathbf{s}) \;\bullet\; \mathbf{\delta}_{2}(\mathbf{s},\mathbf{t}) + \ldots + \mathrm{\Sigma_{s,t}} \; \mathrm{log}\; \mathrm{P}_{\mathbf{\theta}}(\mathrm{S}_{L}=\mathbf{t}\;\mid\; \mathrm{S}_{L-1}=\mathbf{s}) \;\bullet\; \mathbf{\delta}_{L}(\mathbf{s},\mathbf{t}) \\ &+ \mathrm{\Sigma_{s}} \; \mathrm{log}\; \mathrm{P}_{\mathbf{\theta}}(\mathrm{X}_{1}=\mathbf{x}_{1}\;\mid\; \mathbf{S}_{1}=\mathbf{s}) \;\bullet\; \mathbf{\phi}_{1}(\mathbf{s}) + \ldots + \mathrm{\Sigma_{s}} \; \mathrm{log}\; \mathrm{P}_{\mathbf{\theta}}(\mathrm{X}_{L}=\mathbf{x}_{L}\;\mid\; \mathrm{S}_{L}=\mathbf{s}) \;\bullet\; \mathbf{\phi}_{L}(\mathbf{s}) \end{aligned}$$

## Pseudo code SQUAREM

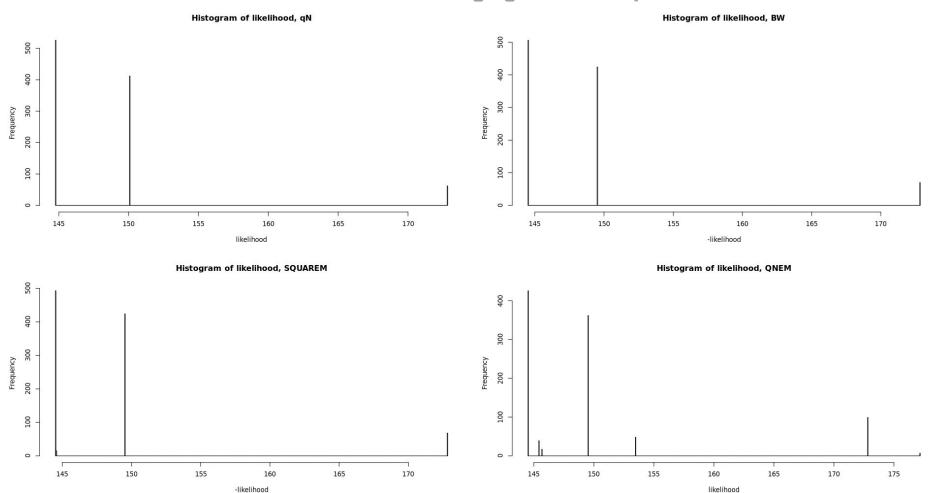
#### Pseudocode:

- 1. Initialise with  $\theta_0$
- 2. 2 EM iterations  $\Rightarrow \theta_1$  et  $\theta_2$
- 3.  $\mathbf{r} = (\theta_0 \theta_1)$
- $4. \quad \mathbf{v} = (\theta_2 \theta_1) \mathbf{r}$
- 5. Compute step length: norm ratio of r and v
- 6. Estimate  $\theta$ ': with  $\theta_0$ , step length, r and v
- 7. Update  $\theta_0$  with EM( $\theta$ ')
- 8. Check convergence
  - a. go back to 1.
  - b. or stop

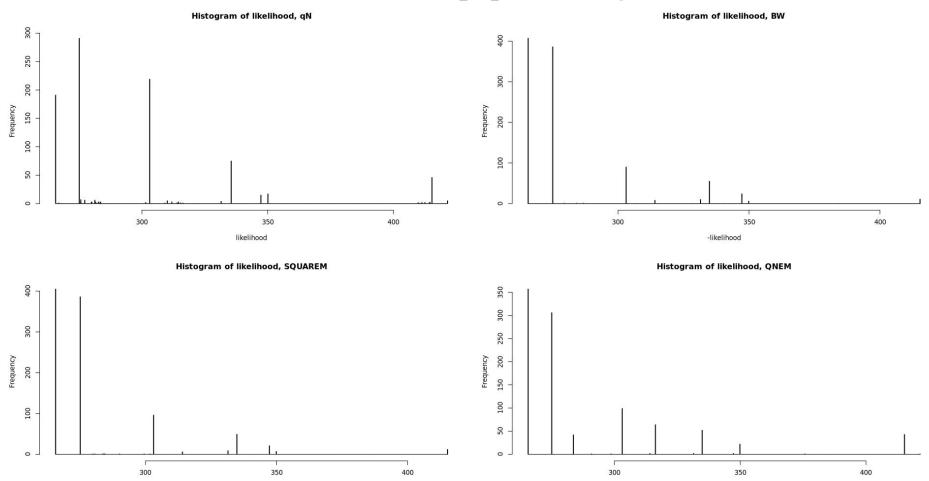
# Umbrella example



# Dichotomised geyser example



# Continuous geyser example



-likelihood

likelihood